

Dev and test set should come from same distribution.

Because dev set is what we use to evaluate model

overfitting  $\rightarrow$  high variance

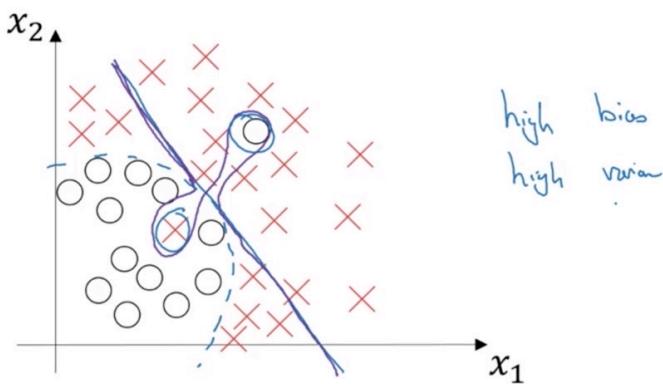
underfitting  $\rightarrow$  high bias

bias ↘

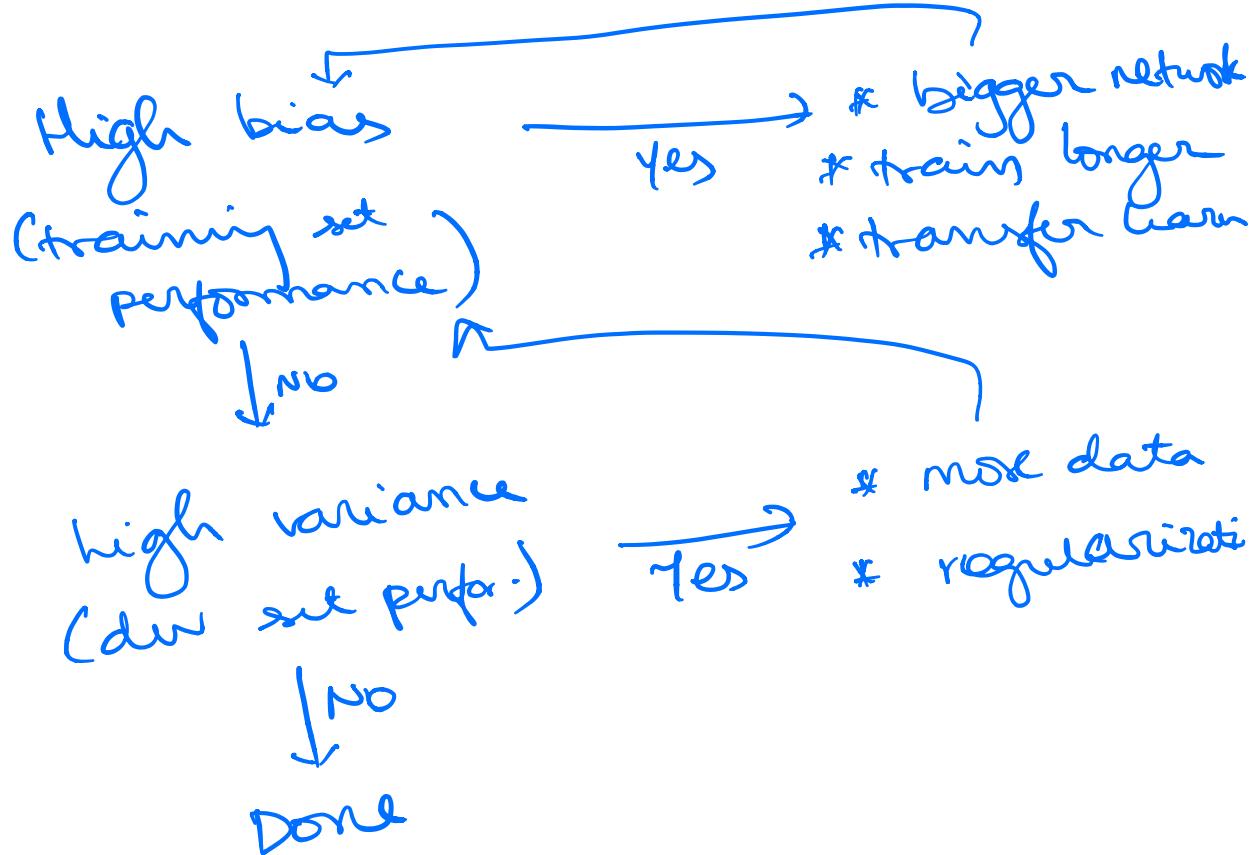
Train set error:	1%	<u>15%</u>	15%	0.5%
Dev set error:	11%	16%	30%	1%
	high variance Hmax: ~10%	high bias	high bias & high varian	low bias low variance

variance ↑

Optimal (Bayes) error: ~ 0%



## Basic recipe:



we don't need to worry about bias - variance trade off as in DL, we can reduce one without affecting the other unlike algs in pre-DL era.

## Regularization:

- \* prevent overfitting
- \* reduce variance

$$J(w, b) = \frac{1}{m} \sum L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2 + \frac{\lambda}{2m} b^2$$

regularization parameter

omit

$L_2$  regularization:

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

$L_1$  regularization:  $\frac{\lambda}{2m} \sum_{i=1}^{n_x} |w_i|$

$w$  will be sparse  
(lots of zeros)

$$= \frac{\lambda}{2m} \|w\|_1$$

↓  
less memory       $\|w\|_n = \left( \sum w_i^n \right)^{1/n}$

$$J(w^{[1]}, b^{[1]}, \dots) = \frac{1}{m} \sum_{i=1}^m L(y_i, \hat{y}_i) + \frac{\gamma}{2m} \sum_{l=1}^L \|w_l\|^2$$

$$\|w\|^2 = \sum_{l=1}^{N_{\text{layer}}} \sum_{j=1}^{n_{l-1}} (w_{lj})^2$$

↳ frobenius norm

for backprop:

$$d w^{[l]} = (\text{unreal}) + \frac{\gamma}{m} w^{[l]}$$

$$w^{[l]} = w^{[l]} - \alpha d w^{[l]}$$

$$= w^{[l]} - \alpha \left( \text{unreal} + \frac{\gamma}{m} w^{[l]} \right)$$

$$= w^{[l]} - \underbrace{\frac{\alpha \gamma}{m} w^{[l]}}_{\text{reduces } w^{[l]}} - \alpha (\text{unreal})$$

$\therefore \gamma$  also called  
weight decay

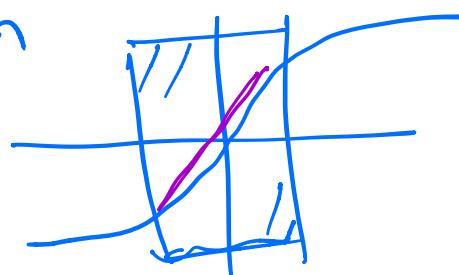
why regularization helps?

Since we have  $\|w\|_F$  added  
(for representing norm)  
to  $J(\theta)$  and we try to  
reduce value of  $J$ , this  
will try to reduce the  
values of  $w$ .

$\therefore$  if  $\lambda$  is large  $\Rightarrow w \xrightarrow{[1]} 0$   
Thus reducing  $w$  values and  
hence we get simpler network

another intuition:

Take tanh fn



If  $\gamma \uparrow \Rightarrow w^{[L3]} \downarrow$   
 $\therefore z^{[L3]} = w^{[L3]} A^{[L-1]} + b^{[L3]}$   
 $\Rightarrow z^{[L3]} \text{ is small}$

$\therefore$  values of  $z^{[L3]}$  will occupy  
 small region of tank.  
 which is almost linear.  
 This will prevent overfitting

Diff between L1 and L2:

- 1) L1 makes model sparse  
 $\rightarrow$  model compression  
 $\rightarrow$  feature selection as coefficients become 0.
- 2) consider a case :  $X : x_1, x_2$   
 and  $x_1$  and  $x_2$  are same  
 so :  $y = x_1$  or  $y = x_2$  is answer.

more generally :  $y = ax_1 + bx_2$   
where  $a+b=1$

$L_1$  gives  $y=x_1$  or  $y=x_2$

$L_2$  gives  $y=0.5x_1 + 0.5x_2$

So, when in real world, due  
to some system error we get  
only  $x_1$  or  $x_2$ , model trained  
with  $L_1$  will suffer while  
model trained with  $L_2$  won't.

$\therefore L_2$  models are more  
robust than  $L_1$ .

Dropout regularization:

with some probability,  
remove some neuron

Implemented using "inverted dropout"

Implementing dropout ("Inverted dropout")

Illustrate with layer  $l=3$ . keep-prob =  $\frac{0.8}{x}$  0.2

$\rightarrow d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep-prob}$

$a_3' = \text{np.multiply}(a_3, d_3)$  #  $a_3 \neq d_3$ .

$\rightarrow a_3' /= \cancel{\text{keep-prob}}$  ← maintains expected value of  $a_3'$   
50 units. ↳ 10 units shut off as same

$$z^{[4]} = w^{[4]} \cdot \frac{a^{[3]}}{l} + b^{[4]}$$

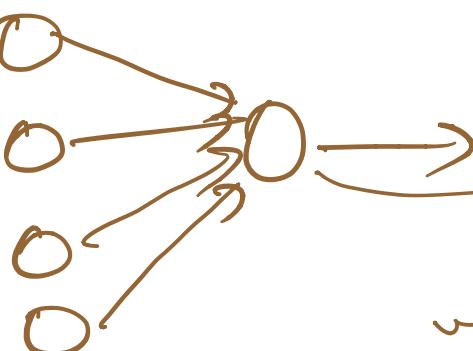
↑ reduced by 20%. Test

$$l = \underline{0.8}$$

no dropout when testing

Downside:

$z$  is not  
well defined



Since nodes are shut off randomly, the node can't place all bets on one feature. So weights become shared

reduce overfitting :

- \* data augmentation
- \* early stopping

Normalizing:

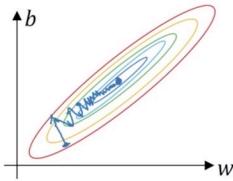
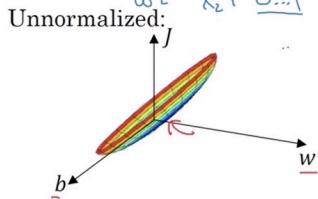
Subtract mean  $\mu = \frac{1}{m} \sum x^{(i)}$   
normalize variance

$$\sigma^2 = \frac{1}{m} \sum (x^{(i)} - \mu)^2$$

$$x = \frac{x - \mu}{\sigma}$$

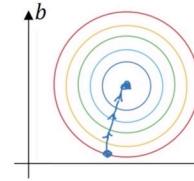
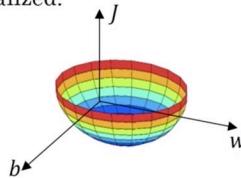
Why normalize inputs?

$$\begin{array}{ll} w_1 & x_1: \underline{1 \dots 1000} \\ w_2 & x_2: \underline{0 \dots 1} \end{array}$$



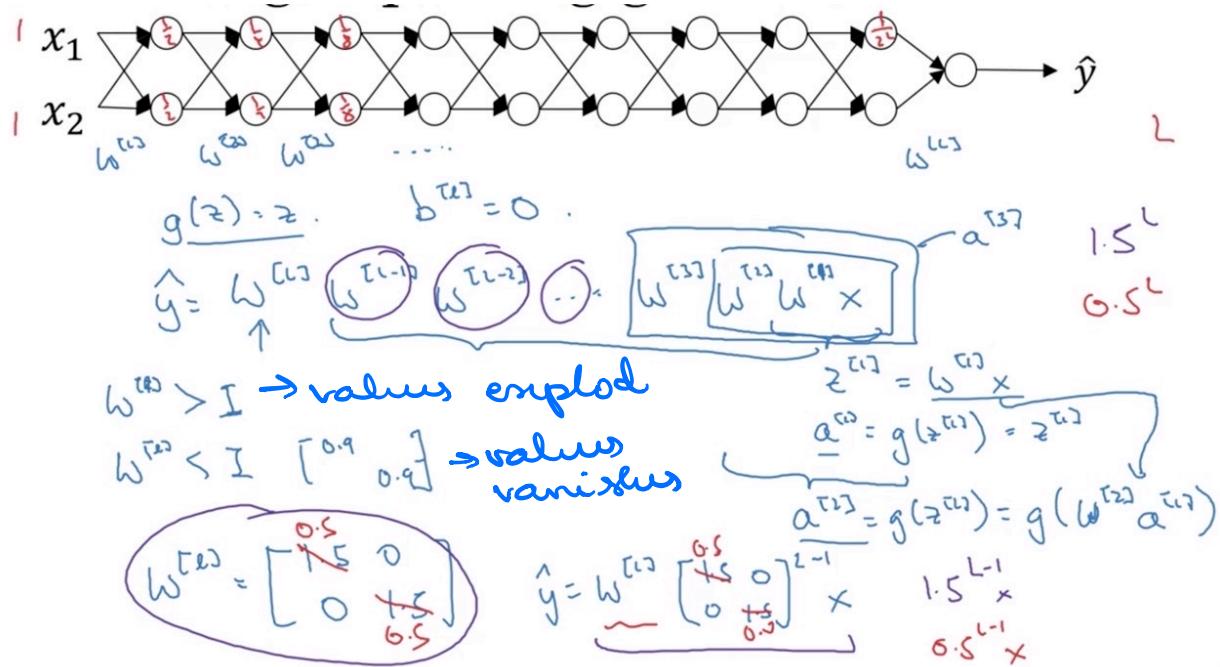
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Normalized:



Andrew Ng

## Vanishing / exploding gradients



## Initialization:

$$\text{variance} = \frac{x}{n}$$

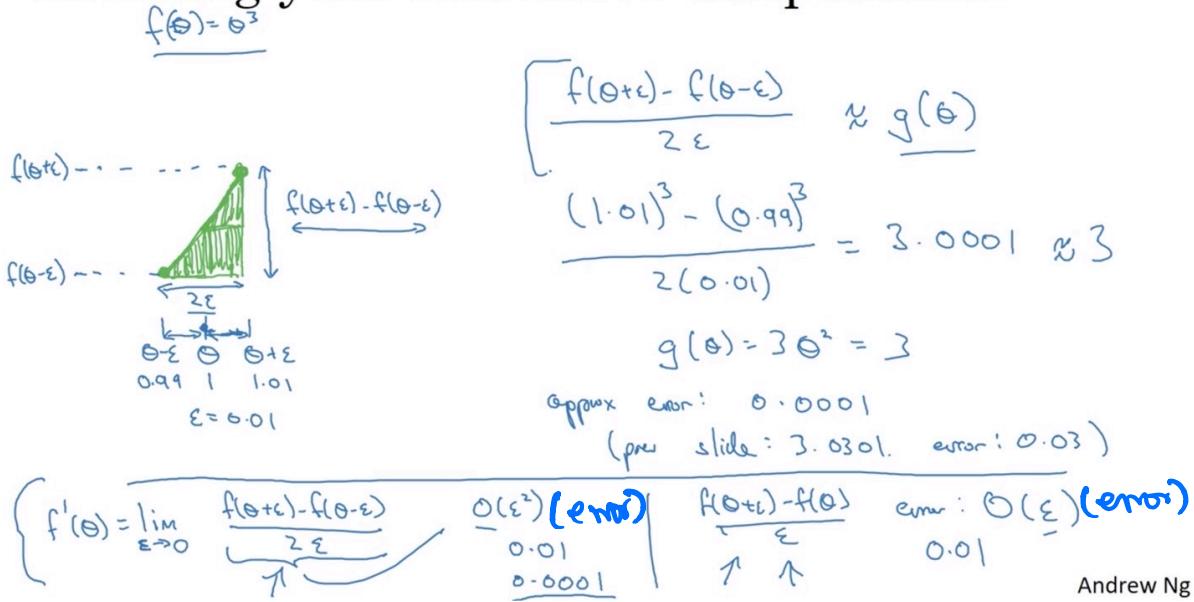
$$w^{[l]} = \text{np.random.rand}() * \text{np.sqrt}\left(\frac{x}{n^{[l-1]}}\right)$$

value :  $x=2$

usually :  $x=1$   
(tanh)

$$\text{Xavier} : \sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

## Checking your derivative computation



## Gradient checking (Grad check)

$$J(\theta) = J(\theta_0, \theta_1, \dots)$$

for each  $i$ :

$$\rightarrow \underline{d\theta_{\text{approx}}[i]} = \frac{J(\theta_0, \theta_1, \dots, \theta_i + \epsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\approx \underline{d\theta[i]} = \frac{\partial J}{\partial \theta_i}$$

$\vdots$

$d\theta_{\text{approx}} \approx d\theta$

Check

$$\rightarrow \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$$

$$\epsilon = 10^{-7}$$

$\times$   $10^{-7} - \text{great!}$

$$\rightarrow 10^{-3} - \text{worry.}$$

Mini batch gradient descent:

split training data into  
small batches

$$X = \{ \underset{\substack{\uparrow \\ \text{one mini batch}}}{x^{13}}, x^{23}, \dots \}$$

$$\text{minibatch } t : \underset{\substack{\downarrow \\ (n_x, \text{size of mini batch})}}{x^{t3}} - \underset{\substack{\downarrow \\ (1, \text{size of mini batch})}}{y^{t3}}$$

for  $t = 1 \dots \text{no. of mini batches}:$

{ implement one step of gradient  
descent as if our training  
set is  $x^{t3}, y^{t3}$

forward prop on  $x^{t3}$

compute cost  $J^{t3}$

back prop to compute gradient  
wrt  $J^{t3}, x^{t3}, y^{t3}$

update  $w, b$

If mini batch size = m:

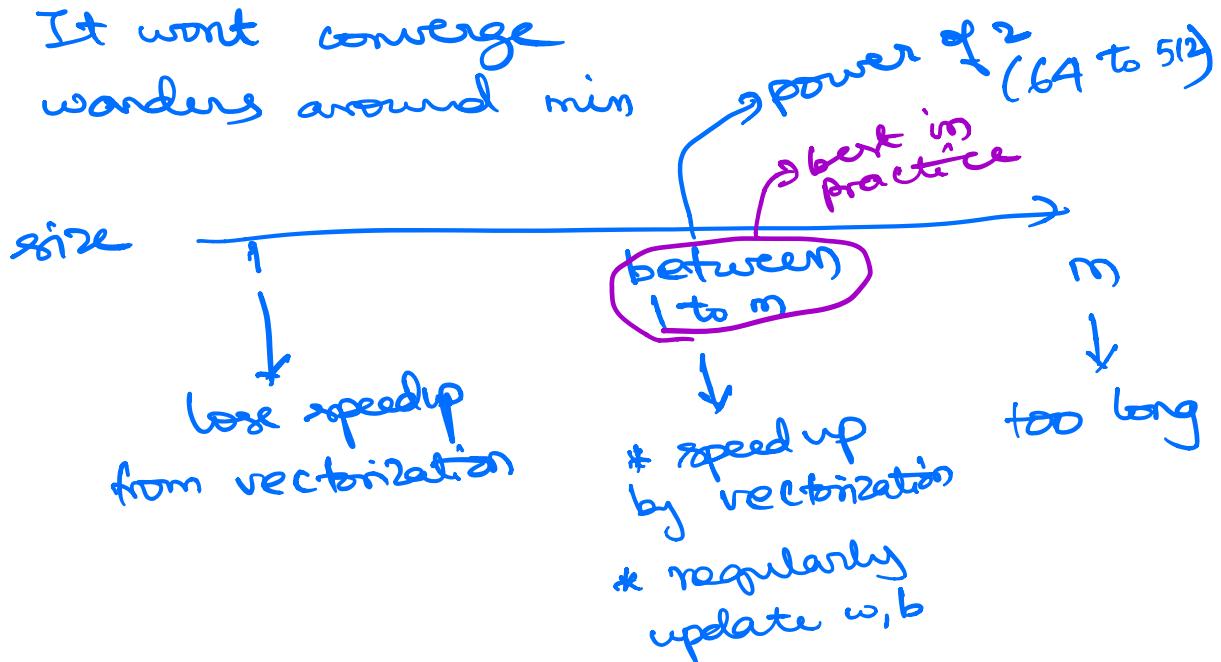
batch gradient descent

else m bs = 1:

stochastic gradient descent

(every example is a mini batch)

It wont converge  
wanders around min



Exponentially weighted moving average:

$$V_t = \beta V_{t-1} + (1-\beta) D_t$$

This is similar to approximately  
averaging over  $\frac{1}{1-\beta}$  day's temp.

$$\therefore \beta = 0.9 \Rightarrow \frac{1}{1-\beta} = \frac{1}{0.1} = 10 \text{ days}$$

$$\beta = 0.98 \Rightarrow \frac{1}{1-\beta} = \frac{1}{0.02} = 50 \text{ days}$$

$$\beta = 0.5 \Rightarrow \frac{1}{1-\beta} = \frac{1}{0.5} = 2 \text{ days}$$

$$V_\theta = 0$$

repeat:

$$\text{get } \theta_t$$

$$V_\theta = \beta V_\theta + (1-\beta)\theta_t$$

} just store one array in memory

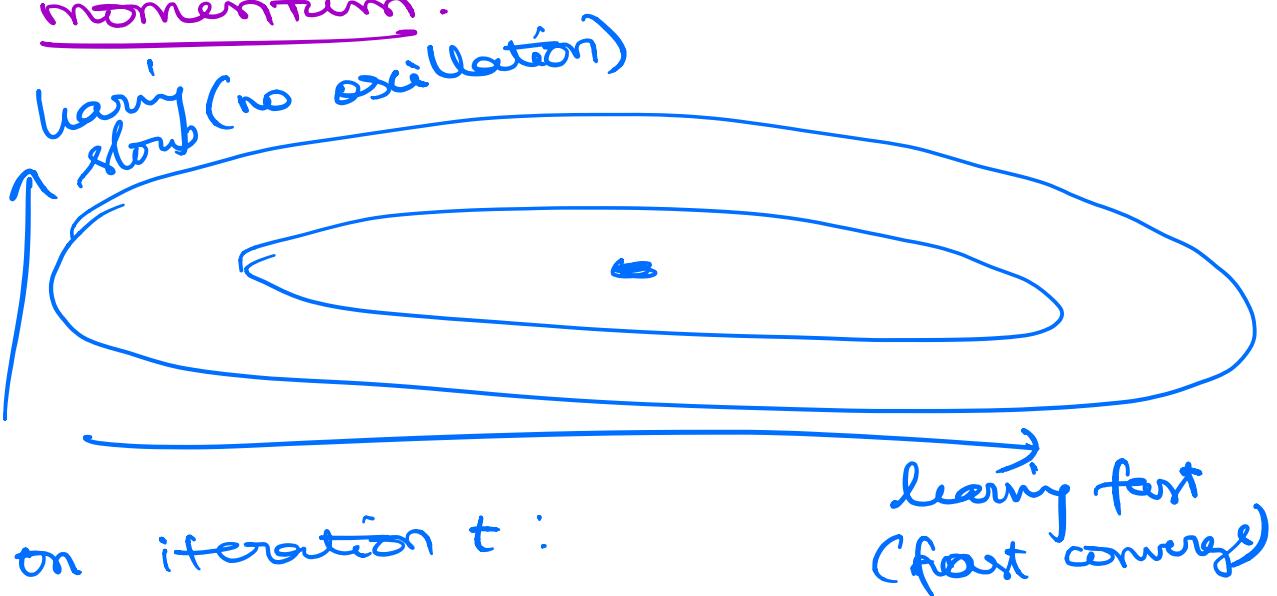
Bias correction:

since  $V_\theta$  is set to 0 initially,  
the starting part of the curve  
will be somewhat underfitting

$$\therefore \text{we can use } v_\theta = \frac{v_\theta}{1 - \beta^t}$$

so, in initial states, the denom will scale the values correctly as  $t$  becomes large, denominator approaches 1 and thus has no effect on  $v_\theta$  values.

Gradient descent with momentum:



on iteration  $t$ :

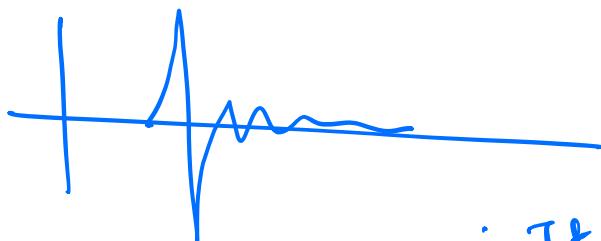
compute  $d_w, d_b$

$$v_d w = \beta v_d w + (1-\beta) d_w$$

$$v_d b = \beta v_d b + (1-\beta) d_b$$

$$\omega = \omega - \alpha V d\omega$$

$$b = b - \alpha V b$$



averaging vertical axis will be almost zero

$\therefore$  It takes less oscillations along vertical axis while maintaining faster convergence

usually  
 $\beta = 0.9$

$$V d\omega : \beta V d\omega + (1-\beta) d\omega$$

↑  
friction      ↑ velocity      ↑ acceleration

no need for bias correction as when  $\beta = 0.9$ , it will work correctly after  $\approx 10$  iteration.

## Root mean square prop (RMS prop):

on iteration t :

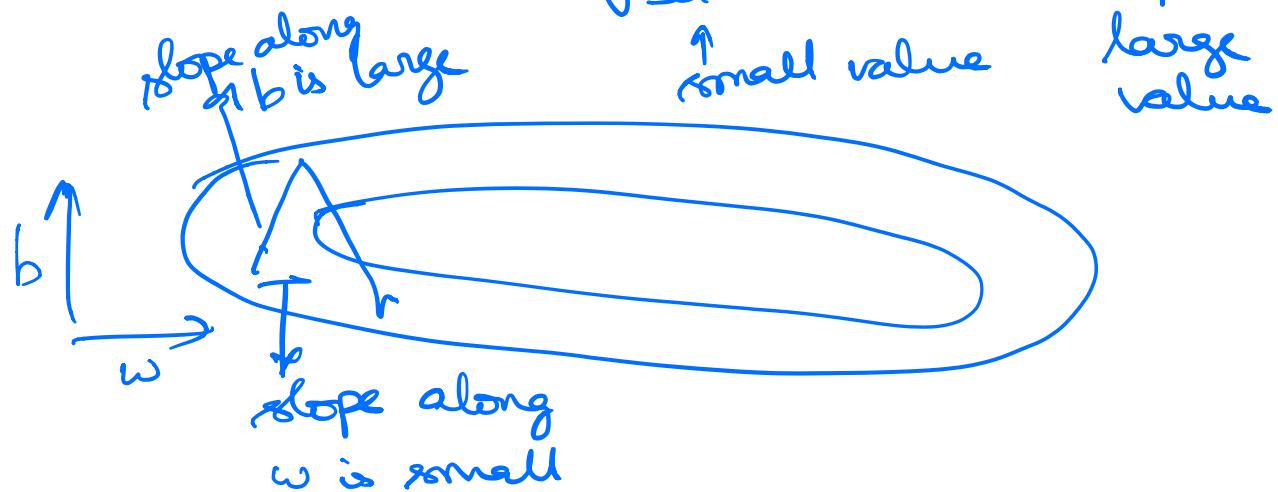
$\epsilon \Rightarrow$  small value  
so that denomi  
is not zero

compute  $d_w, d_b$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) d_w^2$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) d_b^2$$

$$w = w - \frac{\alpha d_w}{\sqrt{S_{dw}} + \epsilon} \quad b = b - \frac{\alpha d_b}{\sqrt{S_{db}} + \epsilon}$$



Adam optimization: (Adaptive moment estimation)

$$V_{dw} = S_{dw} = V_{db} = S_{db} = 0$$

on iteration  $t$ :

compute  $dw, db$

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

$$\begin{matrix} \beta_1 &= 0.9 \\ \beta_2 &= 0.99 \\ \alpha &= 0.001 \end{matrix}$$

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{1 - \beta_1^t}$$

$$V_{db}^{\text{corrected}} = \frac{V_{db}}{1 - \beta_1^t}$$

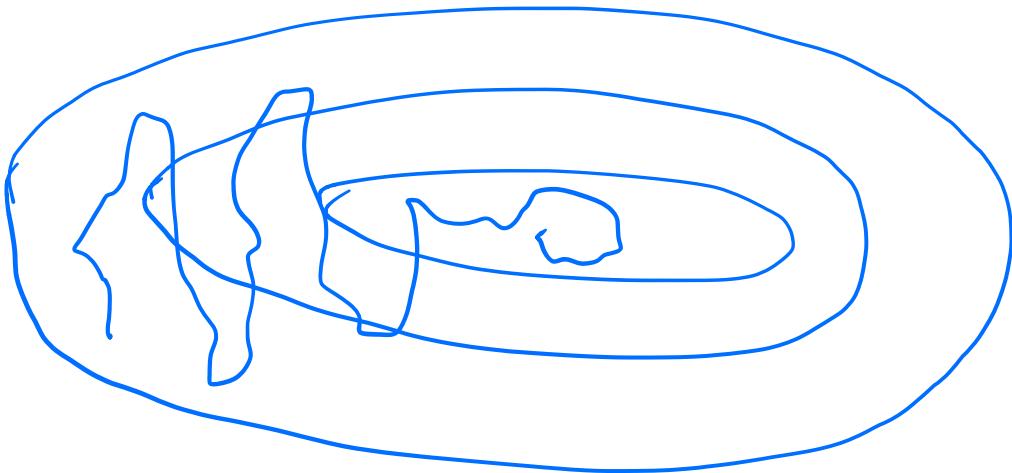
$$S_{dw}^{\text{corrected}} = \frac{S_{dw}}{1 - \beta_2^t}$$

$$S_{db}^{\text{corrected}} = \frac{S_{db}}{1 - \beta_2^t}$$

$$w = w - \frac{\alpha \nabla w^{\text{corrected}}}{\sqrt{S_{ww}^{\text{corrected}}} + \epsilon}$$

$$b = b - \alpha \frac{\nabla b^{\text{corrected}}}{\sqrt{S_{bb}^{\text{corrected}}} + \epsilon}$$

Learning rate decay:



If we slowly reduce  $\alpha$  per iteration:

initial steps: take larger steps

near convergence: take smaller steps

$$\alpha = \frac{1}{1 + \text{decay rate} * \text{epoch num}} * \alpha_0$$

Other ways:

$$\alpha = 0.95^{\text{epoch num}} \cdot \alpha_0 \quad (\text{exponential decay})$$

formula

$$\alpha = \frac{k}{\text{epoch num or } t} \alpha_0$$

discrete decay  $\alpha = \frac{-}{t}$

### Tuning:

1.  $\alpha$
  2.  $\beta$ , minibatch size, # hidden units
  3. # layers, learning rate decay
- 
1. consider random sampling of hyper parameters
  2. consider coarse to fine search.

log scale for  $\alpha$ :



$$= -4$$

sample  $r \in [a, b]$

$$\text{and } \alpha = 10^r$$

for  $\beta$ :

scale using  $1-\beta$

$$\beta = 0.9 \dots 0.999$$

$$1-\beta = 0.1 \dots 0.001$$

$$a = -1$$

$$b = -3$$

$$r \in [-3, -1]$$

$$1-\beta = 10^r$$

$$\beta = 1 - 10^r$$

Batch normalizing:  $\mu = \frac{1}{m} \sum_i z^{[i]}$   $\sigma^2 = \frac{1}{m} \sum_i (z^{[i]} - \mu)^2$

→ makes model robust  
→ has slight regularization effect

$$z_{\text{norm}}^{[i]} = \frac{z^{[i]} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$z_{\text{norm}}^{[i]} = \frac{z^{[i]} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

This will make mean 0 and var 1

if we want something different

use if for further calculations

$$z^{[i]} = \gamma z_{\text{norm}}^{[i]} + \beta$$

used GD to update  $\gamma$  and  $\beta$  as we do for  $w$  &  $b$

if:

$$\gamma = \sqrt{\sigma^2 + \epsilon} \quad \beta = \mu$$

$$\Rightarrow z^{[i]} = z$$

learnable parameters of the model size:  $(n^{[i]}, 1)$

covariate shift:

$$x \rightarrow y$$

↓  
if  $x$  is distribution changes,  
retrain algo

while testing:

we won't have  $\mu, \sigma$  as we'll  
work on one sample.

so use exponentially weighted  
average of  $\mu$  for all minibatches  
and use it when testing

Softmax layer: (multiclass classi-  
fication)

activation fn:

$$t = e^{z^{[L]}}$$

(Softmax regression)

$$\underset{(c,i)}{\underset{\uparrow}{\sum}} y^{\text{data}} = a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^C t_i} \left( \text{or } \frac{t_i}{\sum t_i} \right)$$

$$\text{loss } f_i : - \sum_{j=1}^C y_j \log \hat{y}_j$$

$$J : \frac{1}{m} \sum \text{loss}$$

Back prop:

Initialize  $d_2^{[L]} = \hat{y} - y$  and continue  
back prop.

# Code example

```
import numpy as np
```

```
import tensorflow as tf
```

```
coefficients = np.array([[1], [-20], [25]])
```

```
w = tf.Variable([0], dtype=tf.float32)
```

```
x = tf.placeholder(tf.float32, [3, 1])
```

```
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0] # (w-5)**2
```

```
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
```

```
init = tf.global_variables_initializer()
```

```
session = tf.Session()
```

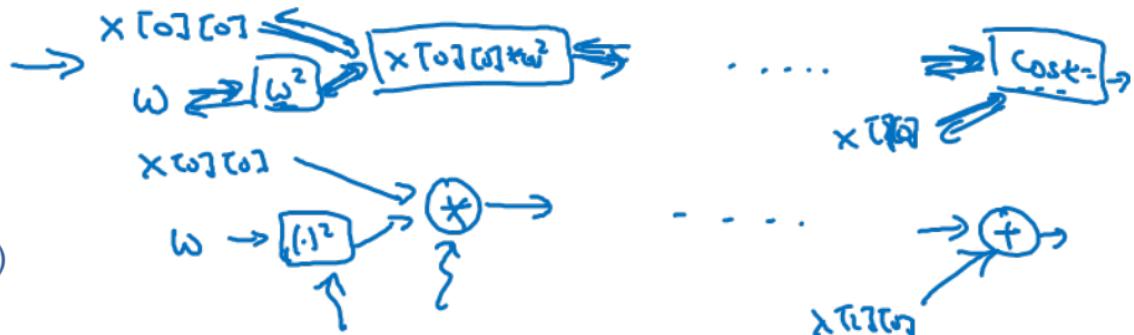
```
session.run(init)
```

```
print(session.run(w))
```

```
for i in range(1000):
```

```
    session.run(train, feed_dict={x: coefficients})
```

```
print(session.run(w))
```



```
with tf.Session() as session:
```

```
    session.run(init)
```

```
    print(session.run(w))
```